



**COMP9242**  
**Advanced Operating Systems**  
**S2/2014 Week 7:**  
**Microkernel Design**



**Australian Government**

COMP9242 S2/2014



# Copyright Notice

---



## These slides are distributed under the Creative Commons Attribution 3.0 License

- You are free:
  - to share—to copy, distribute and transmit the work
  - to remix—to adapt the work
- under the following conditions:
  - **Attribution:** You must attribute the work (but not in any way that suggests that the author endorses you or your use of the work) as follows:
    - “Courtesy of Gernot Heiser, [Institution]”, where [Institution] is one of “UNSW” or “NICTA”

The complete license text can be found at  
<http://creativecommons.org/licenses/by/3.0/legalcode>

# Microkernel Principles: Minimality



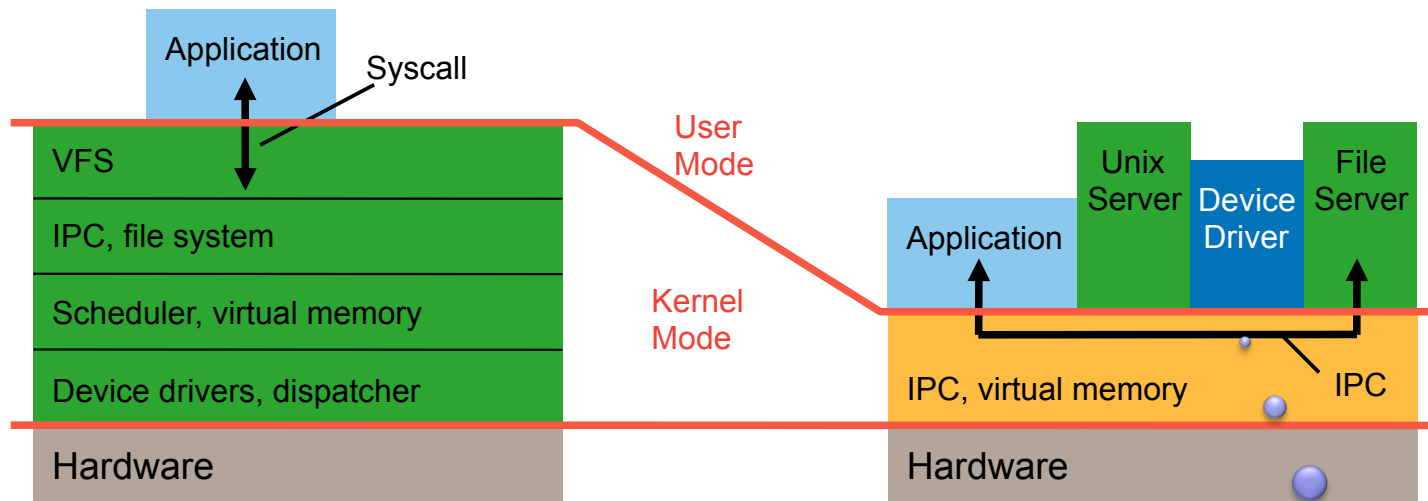
*A concept is tolerated inside the microkernel only if moving it outside the kernel, i.e. permitting competing implementations, would prevent the implementation of the system's required functionality. [SOSP'95]*

- Advantages of resulting small kernel:
  - Easy to implement, port?
  - Easier to optimise
  - Hopefully enables a minimal *trusted computing base* (TCB)
  - Easier debug, maybe even *prove* correct?
- Challenges:
  - API design: generality despite small code base
  - Kernel design and implementation for high performance

Limited by arch-specific micro-optimisations

Small attack surface, fewer failure modes

# Consequence of Minimality: User-level Services



IPC performance is critical!

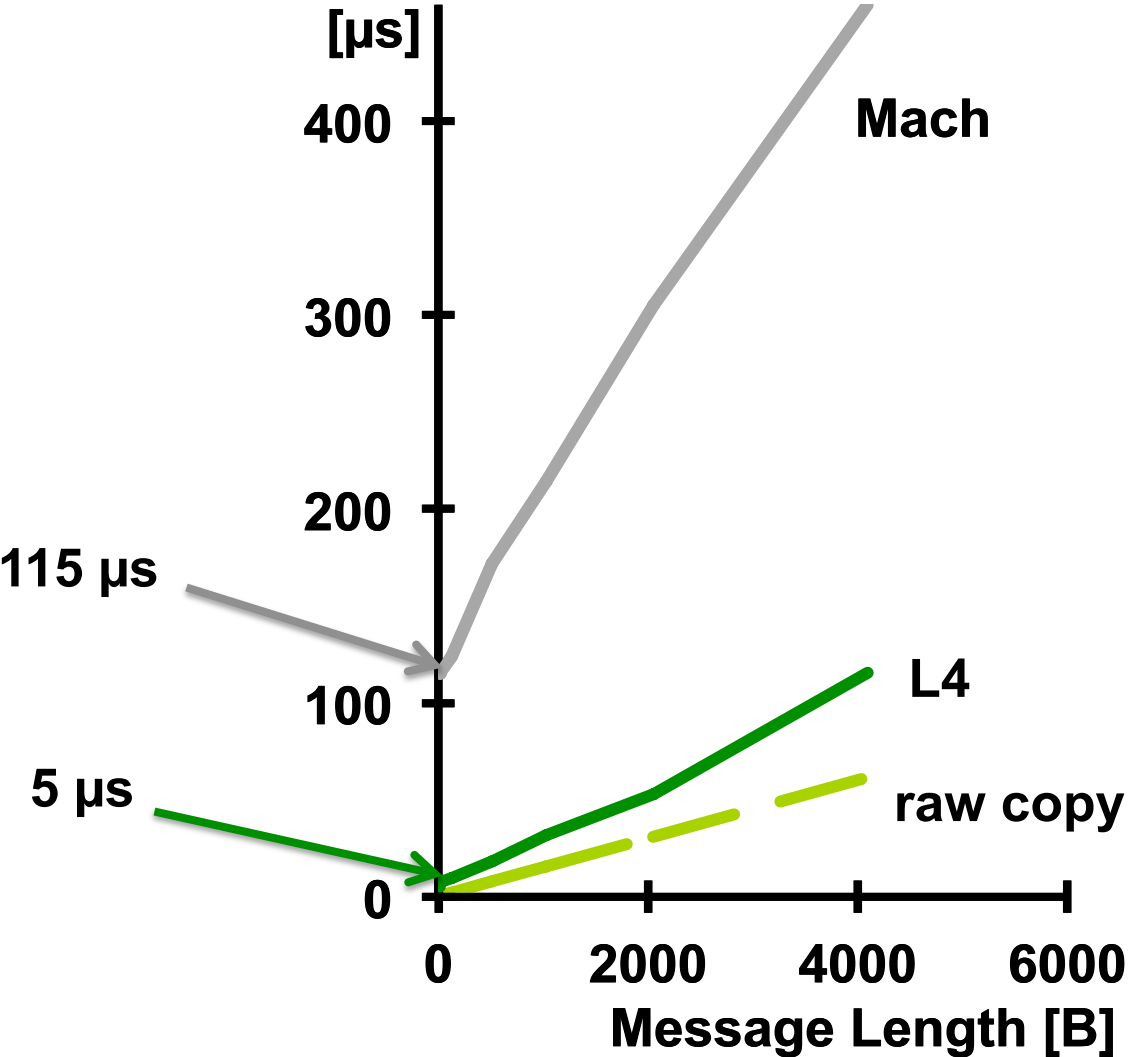
- Kernel provides no services, only mechanisms
- Kernel is policy-free
  - Policies limit (good for 90% of cases, disastrous for some)
  - “General” policies lead to bloat, inefficiency

# 1993 “Microkernel” IPC Performance



i486 @  
50 MHz

Culprit:  
Cache  
footprint  
[SOSP'95]



# L4 IPC Performance over 20 Years



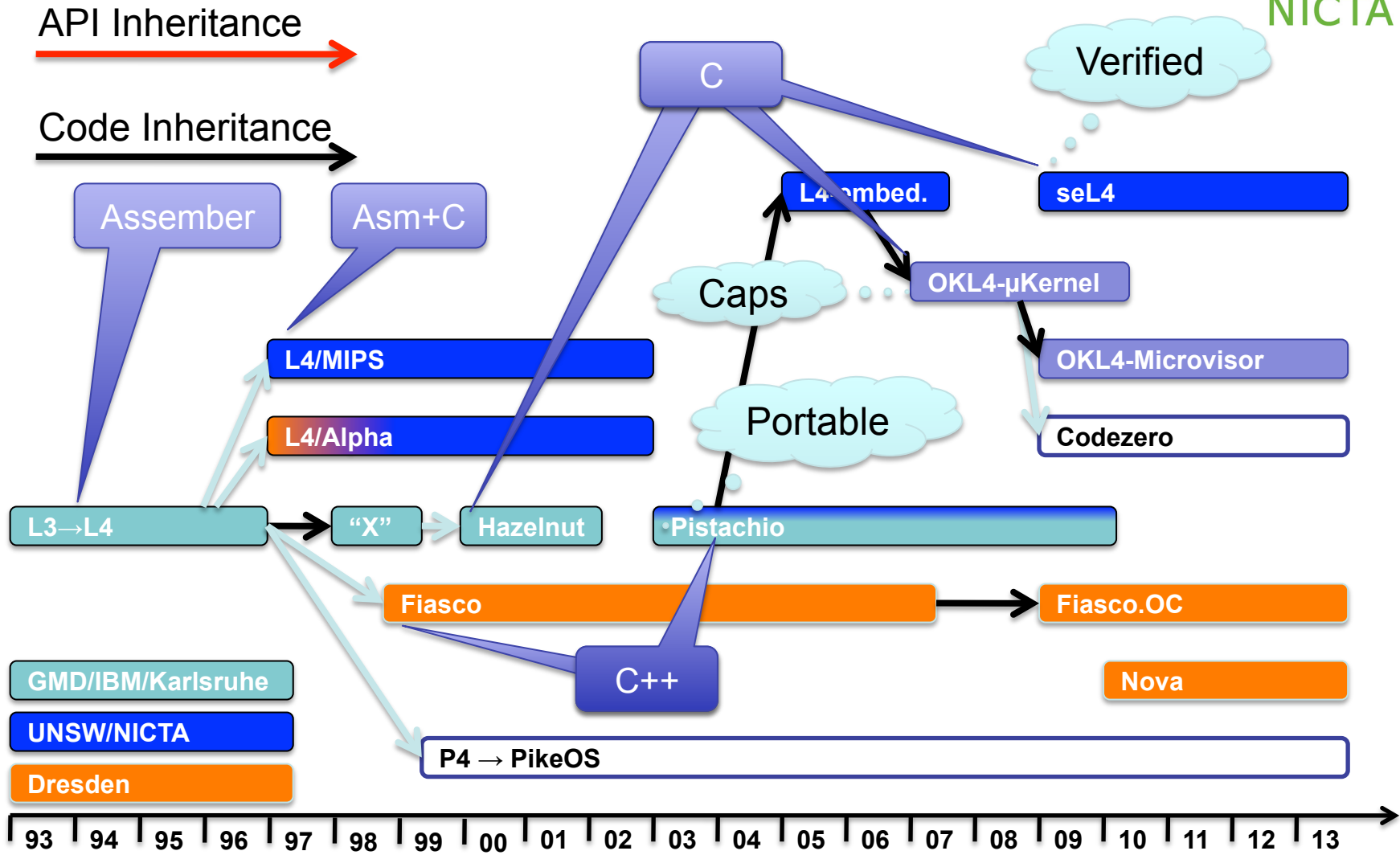
Name	Year	Processor	MHz	Cycles	µs
Original	1993	i486	50	250	5.00
Original	1997	Pentium	160	121	0.75
L4/MIPS	1997	R4700	100	86	0.86
L4/Alpha	1997	21064	433	45	0.10
Hazelnut	2002	Pentium 4	1,400	2,000	1.38
Pistachio	2005	Itanium	1,500	36	0.02
OKL4	2007	XScale 255	400	151	0.64
NOVA	2010	i7 Bloomfield (32-bit)	2,660	288	0.11
seL4	2013	i7 Haswell (32-bit)	3,400	301	0.09
seL4	2013	ARM11	532	188	0.35
seL4	2013	Cortex A9	1,000	316	0.32

# Minimality: Source Code Size



Name	Architecture	C/C++	asm	total kSLOC
Original	i486	0	6.4	6.4
L4/Alpha	Alpha	0	14.2	14.2
L4/MIPS	MIPS64	6.0	4.5	10.5
Hazelnut	x86	10.0	0.8	10.8
Pistachio	x86	22.4	1.4	23.0
L4-embedded	ARMv5	7.6	1.4	9.0
OKL4 3.0	ARMv6	15.0	0.0	15.0
Fiasco.OC	x86	36.2	1.1	37.6
seL4	ARMv6	9.7	0.5	10.2

# L4 Family Tree





# L4 Deployments – in the Billions



## SiMKo 3 “Merkelphone”



# L4 Design and Implementation



## Implement. Tricks [SOSP'93]

- Process kernel
- Virtual TCB array
- Lazy scheduling
- Direct process switch
- Non-preemptible
- Non-portable
- Non-standard calling convention
- Assembler

## Design Decisions [SOSP'95]

- Synchronous IPC
- Rich message structure, arbitrary out-of-line messages
- Zero-copy register messages
- User-mode page-fault handlers
- Threads as IPC destinations
- IPC timeouts
- Hierarchical IPC control
- User-mode device drivers
- Process hierarchy
- Recursive address-space construction

**Objective: Minimise cache footprint and TLB misses**

# DESIGN

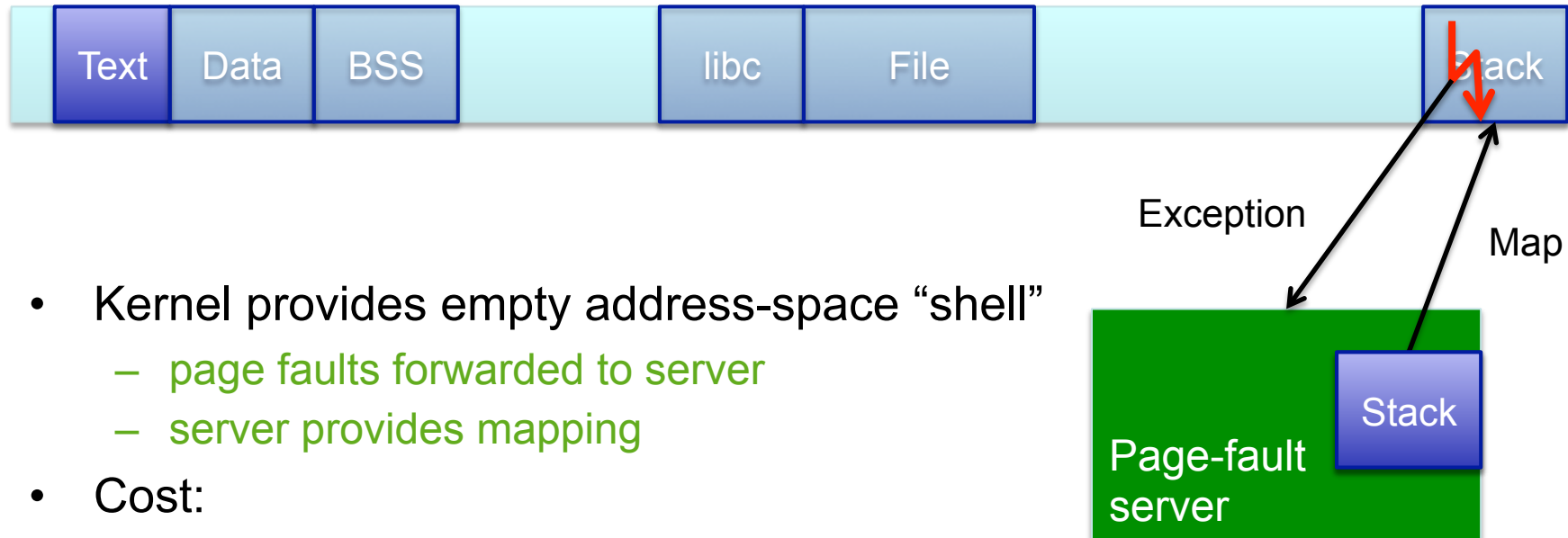
# What Mechanisms?

---



- Fundamentally, the microkernel must abstract
  - *Physical memory*
  - *CPU*
  - *Interrupts/Exceptions*
- Unfettered access to any of these bypasses security
  - No further abstraction needed for devices
    - memory-mapping device registers and interrupt abstraction suffices
    - ...but some generalised memory abstraction needed for I/O space
- Above isolates execution units, hence microkernel must also provide
  - *Communication* (traditionally referred to as *IPC*)
  - *Synchronization*

# Policy-Free Address-Space Management



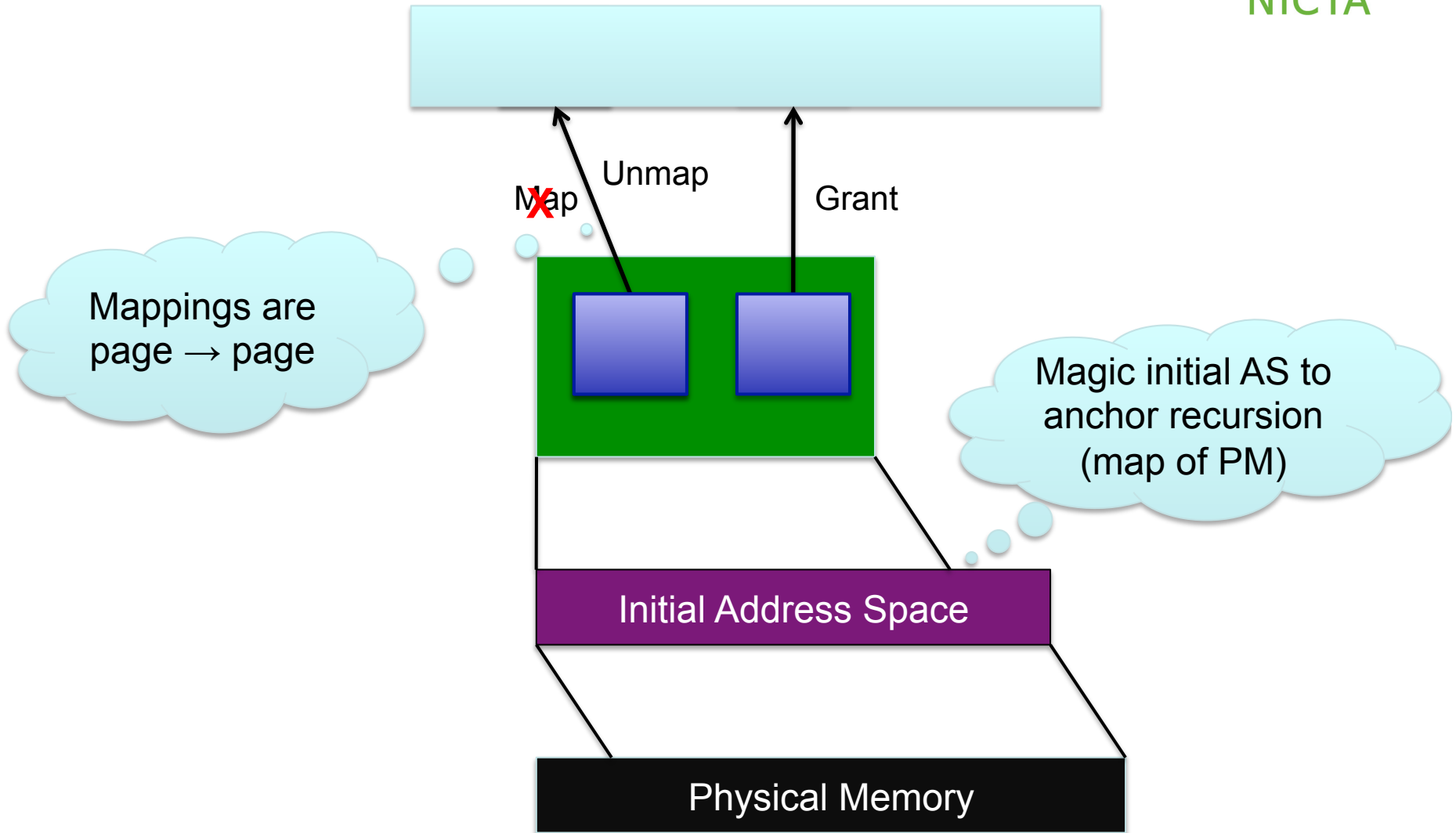
- Kernel provides empty address-space “shell”
  - page faults forwarded to server
  - server provides mapping
- Cost:
  - 1 round-trip IPC, plus mapping operation
    - mapping may be side effect of IPC
    - kernel may expose data structure
  - kernel mechanism for forwarding page-fault exception
- “External pagers” first appeared in Mach [Rashid et al, '88]
  - ... but were optional (and slow) – in L4 there’s no alternative

# Abstracting Memory: Address Spaces



- Minimum address-space abstraction: empty slots for page mappings
  - paging server can fill with mappings
    - virtual address  $\rightarrow$  physical address + permissions
- Can be
  - page-table-like: array under full user control
  - TLB-like: cache for mappings which may vanish
- Main design decision: is source of a mapping a page or a frame?
  - Frame: hardware-like
  - Page: recursive address spaces (original L4 model)

# Traditional L4: Recursive Address Spaces



# Recursive Address Space Experience



## API complexity: Recursive address-space model

- Conceptually elegant
  - trivially supports virtualization
- Drawback: Complex mapping database
  - Kernel needs to track mapping relationships
    - Tear down dependent mappings on unmap
  - Mapping database problems:
    - accounts for 1/4–1/2 of kernel memory use
    - SMP coherence is performance bottleneck
- NICTA's L4-embedded, OKL4 removed MDB
  - Map frames rather than pages
    - need separate abstraction for frames / physical memory
    - subsystems no longer virtualizable (even in OKL4 cap model)
- Properly addressed by seL4's capability-based model
  - But have cap derivation tree, subject of on-going research





# Abstracting Execution



- Can abstract as:
  - kernel-scheduled threads
    - Forces (scheduling) policy into the kernel
  - vCPUs or scheduler activations
    - This essentially virtualizes the timer interrupt through upcall
      - Scheduler activations also upcall for exceptions, blocking etc
    - Multiple vCPUs only for real multiprocessing
- Threads can be tied to address space or “migrating”

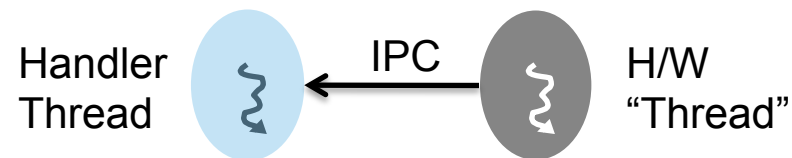
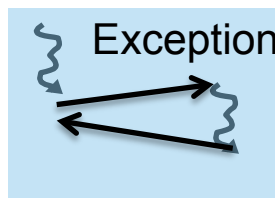


- Implementation-wise not much of a difference
- Tight integration/interdependence with IPC model!

# Abstracting Interrupts and Exceptions



- Can abstract as:
  - Upcall to interrupt/exception handler
    - hardware-like diversion of execution
    - need to save enough state to continue interrupted execution
  - IPC message to handler from magic “hardware thread”
    - OS-like
    - needs separate handler thread ready to receive

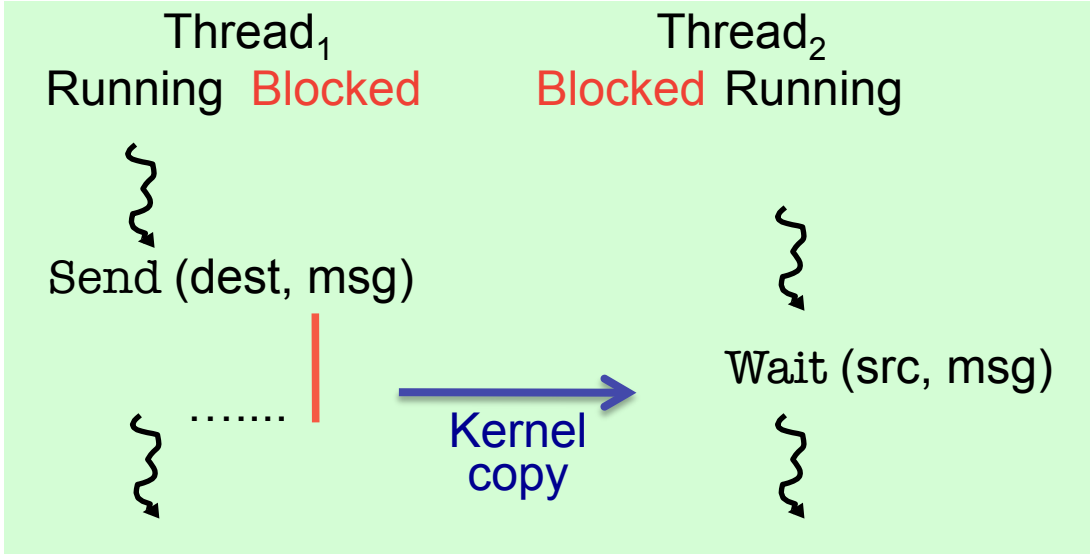


- Page fault tends to be special-cased for practical reason
  - Tends to require handling external to faulter
    - IPC message to page-fault server rather than exception handler
  - But also “self-paging” as in Nemesis [Hand '99] or Barrelfish

# L4 Synchronous IPC



**Rendezvous model**

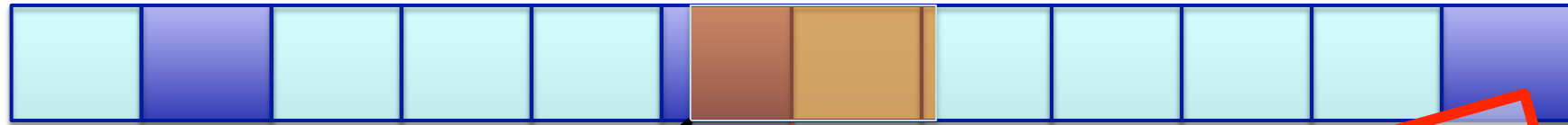


- Kernel executes in sender's context
- copies memory data directly to receiver (single-copy)
- leaves message registers unchanged during context switch (zero copy)

# “Long” IPC



Sender address space



Kernel copy

Receiver address space



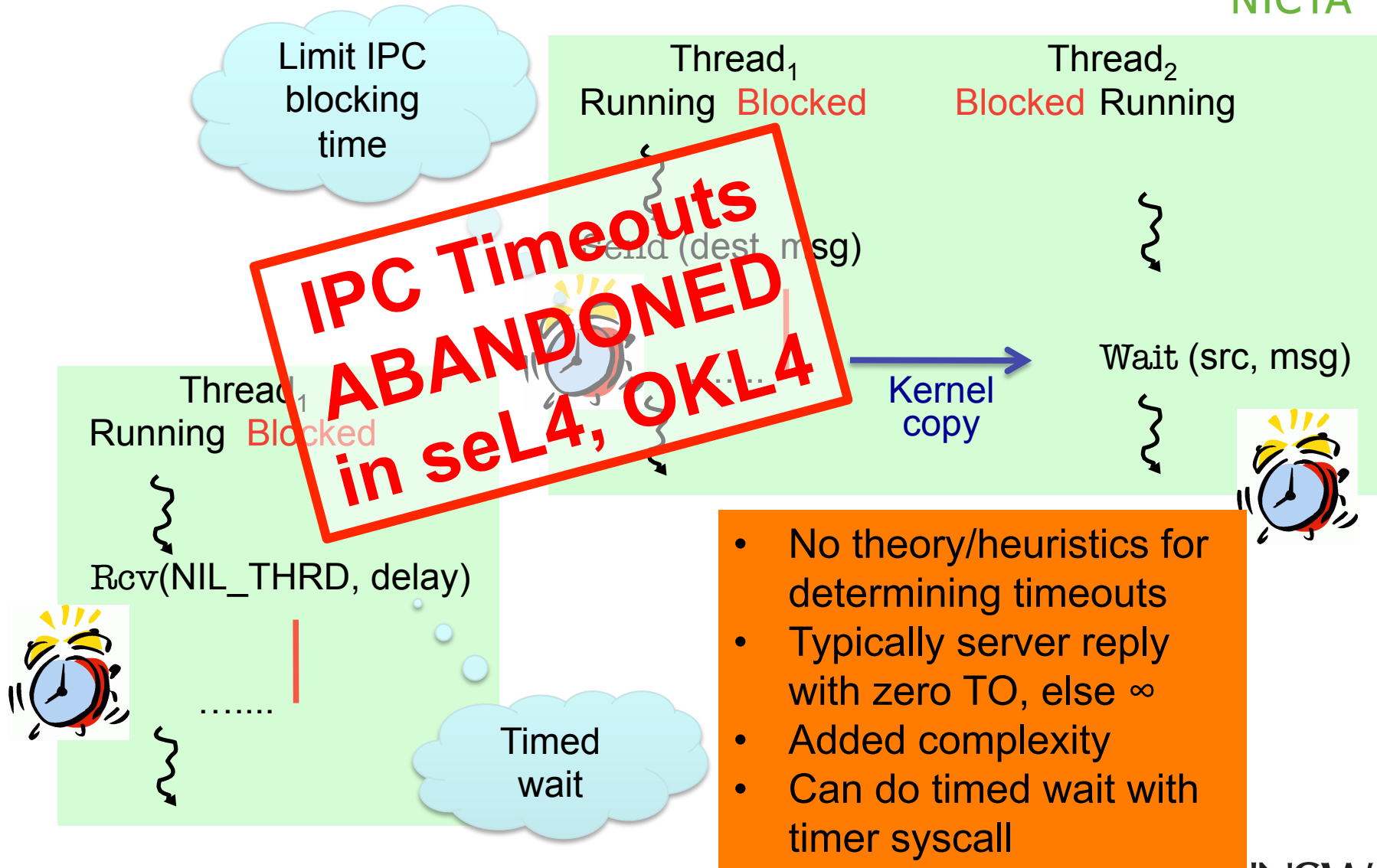
**LONG IPC  
ABANDONED**

- IPC page faults are nested exceptions  $\Rightarrow$  In-kernel concurrency
  - L4 executes with interrupts disabled for performance, no concurrency
- Must invoke untrusted usermode page-fault handlers
  - potential for DOSing other thread
- Timeouts to avoid DOS attacks
  - complexity

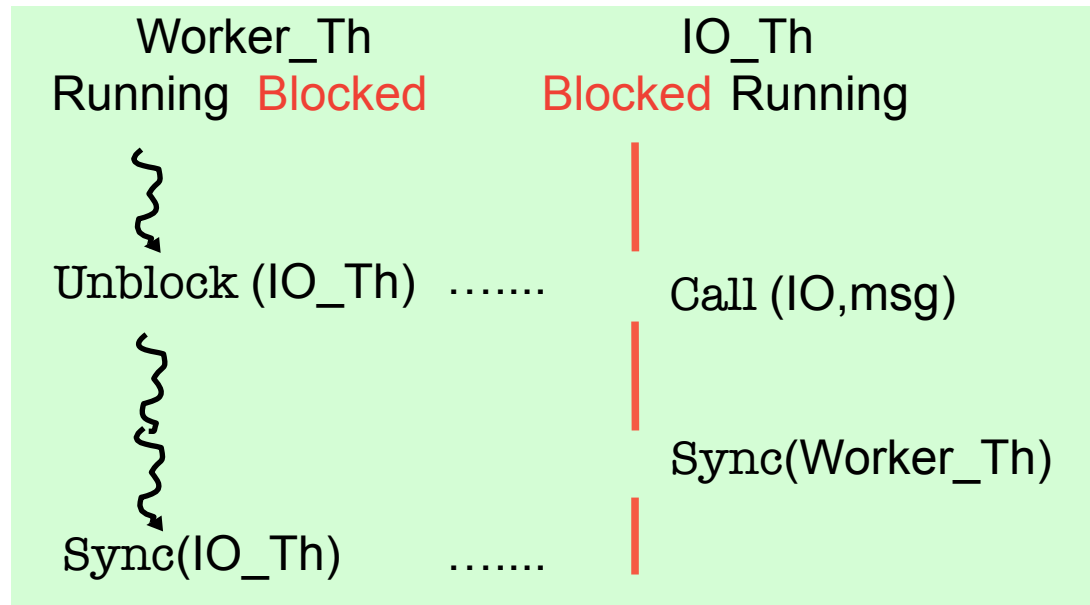
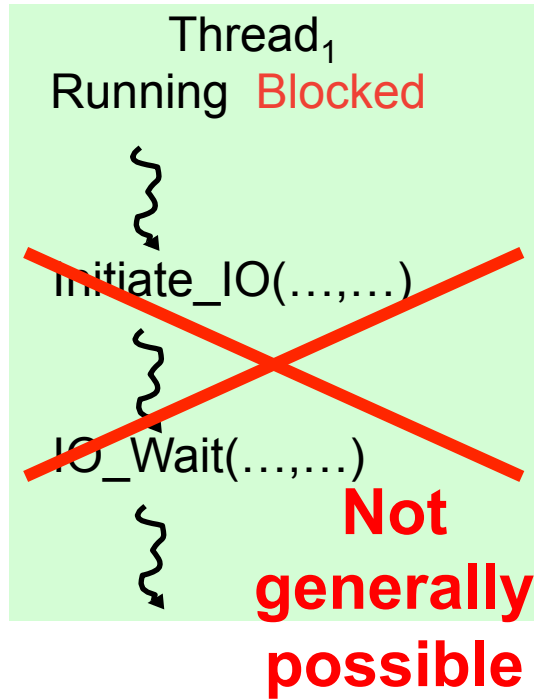
Why have long IPC?

- POSIX-style APIs
  - write (fd, buf, nbytes)
- Usually prefer shared buffers

# Timeouts

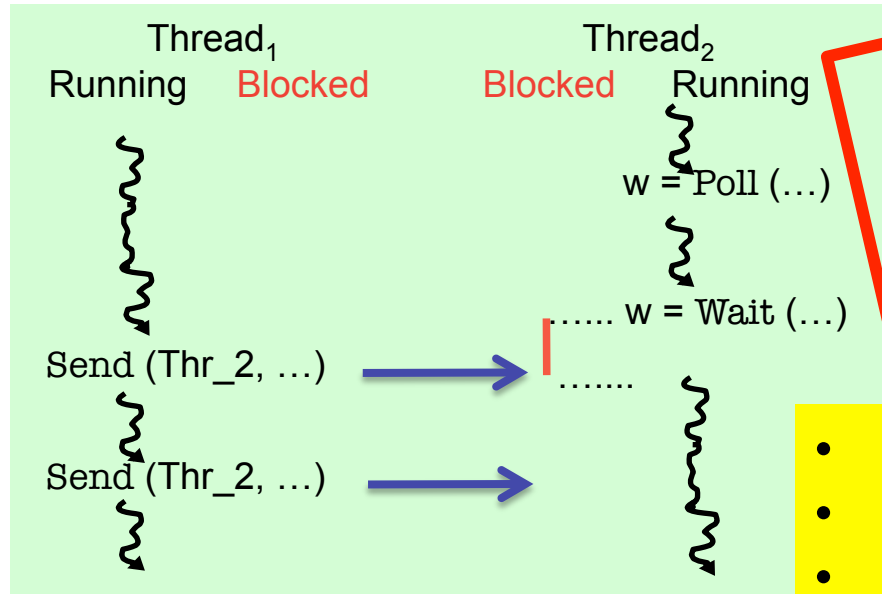


# Synchronous IPC Issues



- Sync IPC forces multi-threaded code!
- Also poor choice for multi-core

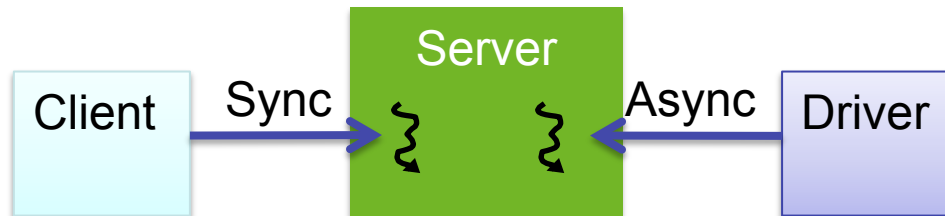
# Asynchronous Notifications



**Sync IPC  
complemented  
with async**

- Delivers few bits (destructively)
- Kernel only buffers single word
- Maps well to interrupts, exceptions

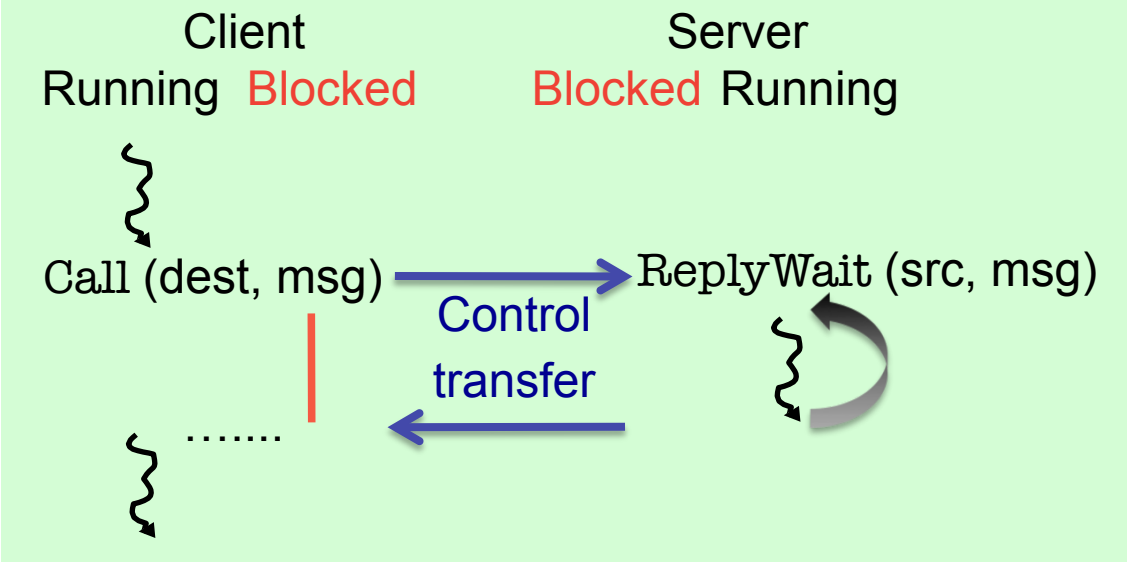
- Thread can wait for synchronous and asynchronous messages concurrently



# Is Synchronous IPC Redundant?



**2 IPC mechanisms:  
Minimality violation**



Sync IPC is useful *intra-core*:

- fast control transfer
- mimics migrating threads
- enables scheduling context donation
  - useful for real-time systems

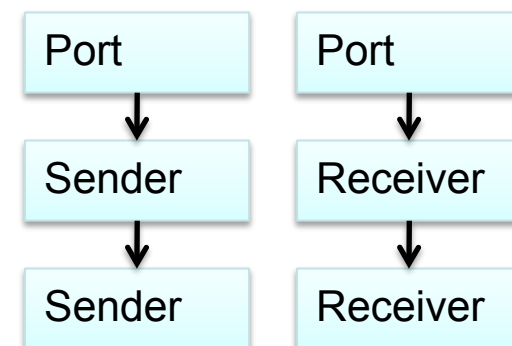
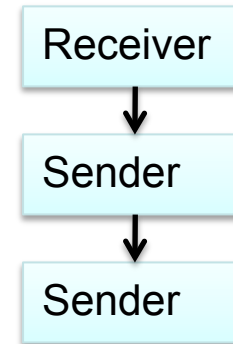
But presently no clean model!



# Direct vs Indirect IPC Addressing



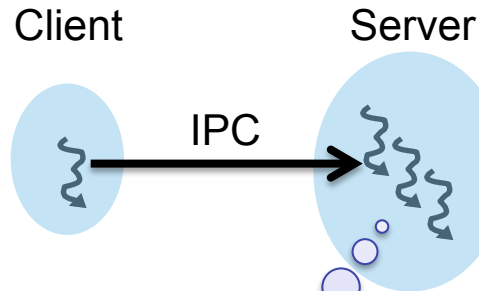
- Direct: Queue senders/messages at receiver
  - Need unique thread IDs
  - Kernel guarantees identity of sender
    - useful for authentication
- Indirect: Mailbox/port object
  - Just a user-level handle for the kernel-level queue
  - Extra object type – extra weight?
  - Communication partners are anonymous
    - Need separate mechanism for authentication



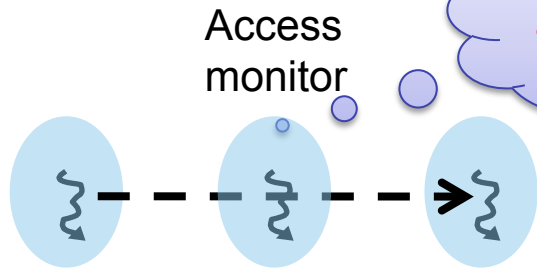
# IPC Destination Naming



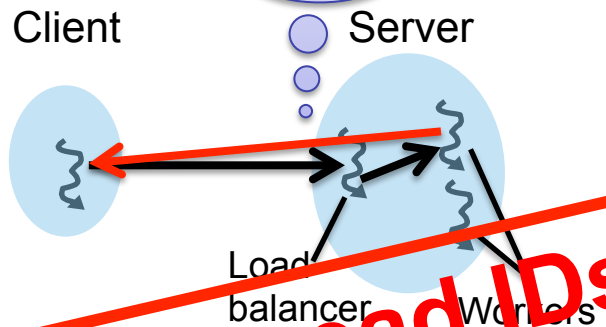
Original L4 addressed IPC to threads



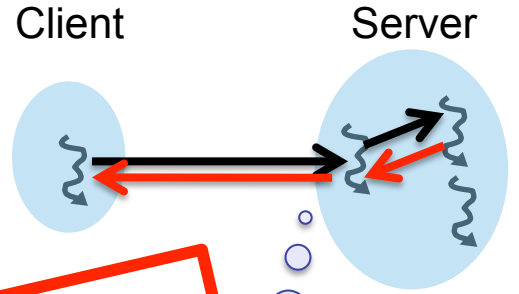
Client must do load balancing?



RPC reply from wrong thread!



**Thread IDs replaced by IPC "endpoints" (ports)**

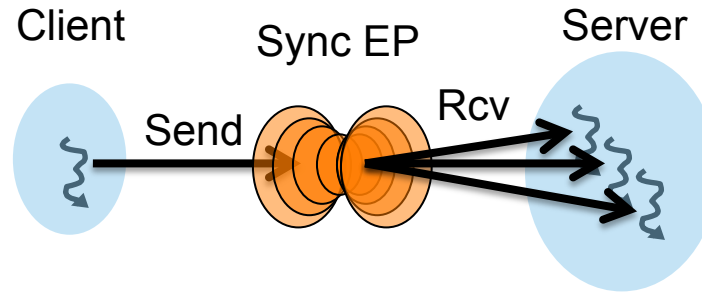
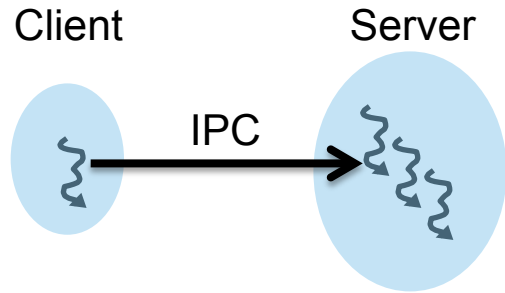


All IPCs duplicated!

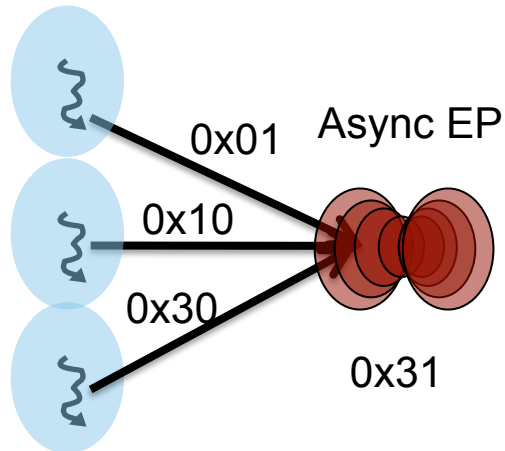
Interpose transparently?

- Inefficient designs
- Poor information hiding
- Covert channels [Shapiro '02]

# Endpoints



- Sync EP queues senders/receivers
- Does not buffer messages



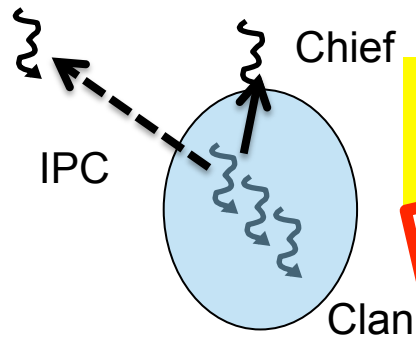
- Async EP accumulates bits

# Other Design Issues



## IPC Control: “Clans & Chiefs”

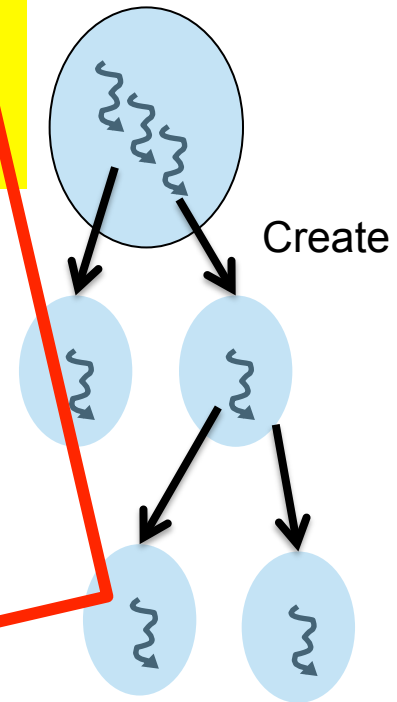
## Process Hierarchy



IPC outside clan  
re-directs to chief

Hierarchical  
resource  
management

**Hierarchies  
replaced by  
delegatable cap-  
based access  
control**

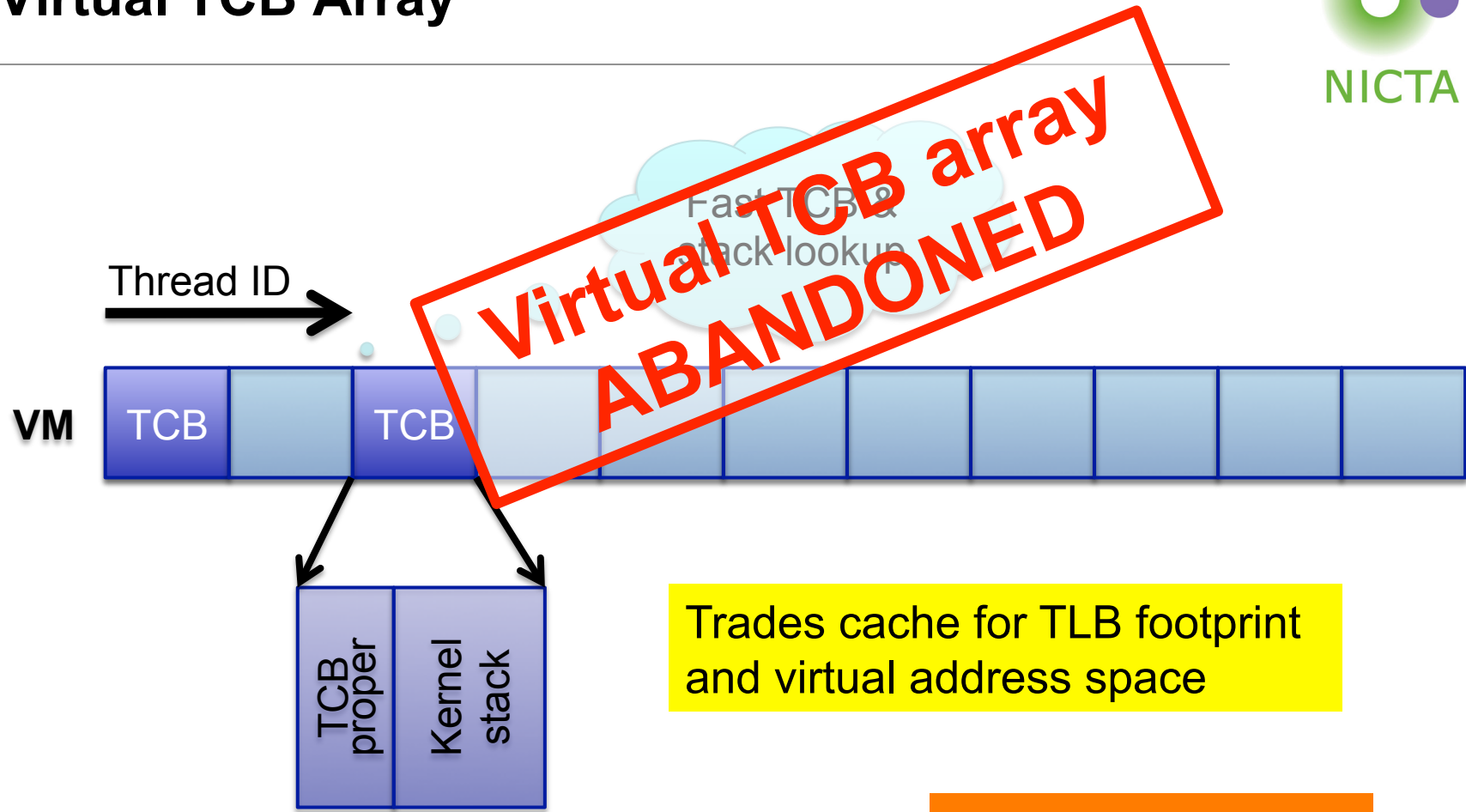


- Inflexible, clumsy, inefficient hierarchies!
- Fundamental problem: no rights delegation

---

# IMPLEMENTATION

# Virtual TCB Array



Trades cache for TLB footprint and virtual address space

Not worthwhile on modern processors!

# Process Kernel: Per-Thread Kernel Stack



Get own TCB base by masking stack pointer

- Not worthwhile on modern processors!
- Stacks can dominate kernel memory use!



- Reduces TLB footprint at cost of cache and kernel memory
- Easier to deal with blocking

# Scheduler Optimisation Tricks: “Lazy Scheduling”

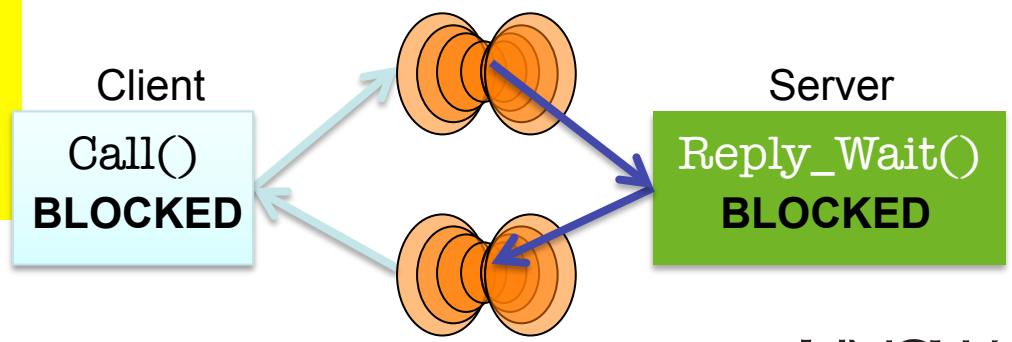


```
thread_t schedule() {  
    foreach (prio in priorities) {  
        foreach (thread in runQueue[prio]) {  
            if (isRunnable(thread))  
                return thread;  
            else  
                schedDequeue(thread);  
        }  
    }  
    return idleThread;  
}
```

Problem: Unbounded scheduler execution time!

Idea: leave blocked threads in ready queue, scheduler cleans up

- Frequent blocking/unblocking in IPC-based systems
- Many ready-queue manipulations





# Alternative: “Benno Scheduling”



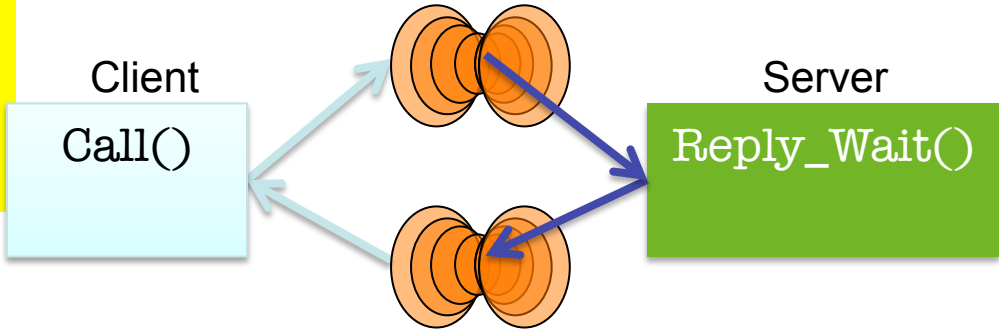
```
thread_t schedule() {  
    foreach (prio in priorities) {  
        foreach (thread in runQueue[prio]) {  
            if (isRunnable(thread))  
                return thread;  
            else  
                schedDequeue(thread);  
        }  
    }  
    return idleThread;  
}
```

**Lazy scheduling REPLACED**

Only current thread needs fixing up at preemption time!

Idea: Lazy on unblocking instead on blocking

- Frequent blocking/unblocking in IPC-based systems
- Many ready-queue manipulations



# Scheduler Optimisation: “Direct Process Switch”



- Sender was running  $\Rightarrow$  had highest prio
- If receiver prio  $\geq$  sender prio  $\Rightarrow$  run receiver

Implication: *Time slice donation* – receiver runs on sender’s time slice

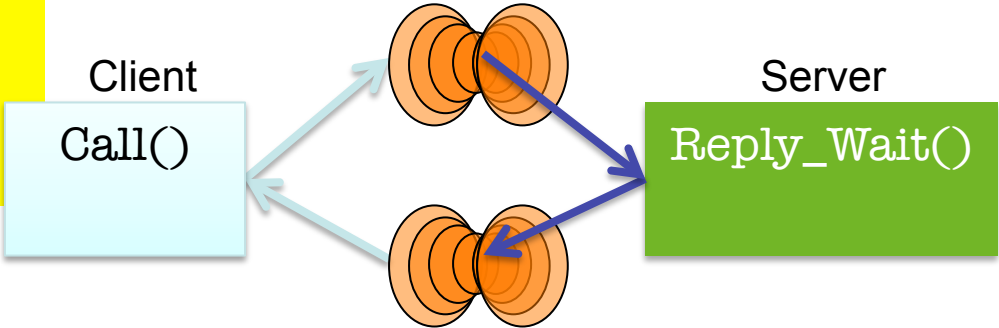
- Arguably, sender should donate back if it’s a server replying to a Call()
- Hence, always donate on Reply\_Wait()

Problem:

- Accounting (RT systems)
- Policy
- No clean model yet!

Idea: Don’t invoke scheduler if you know who’ll be chosen

- Frequent context switches in IPC-based systems
- Many scheduler invocations



# Speaking of Real Time...



- Kernel runs with interrupts disabled
  - No concurrency control ⇒ simpler kernel
    - Easier reasoning about correctness
    - Better average-case performance

- How about long-running system calls?
  - Use strategic *preemption points*
  - (Original) Fiasco has fully preemptible kernel

```
while (!done) {  
    process_stuff();  
    PSW.IRQ_disable=1;  
    PSW.IRQ_disable=0;  
}
```

Limited concurrency in kernel!



Lots of concurrency in kernel!

- Like commercial microkernels (QNX, Green Hills INTEGRITY)

# Incremental Consistency

Good fit to event kernel!



**Avoids concurrency in (single-core) kernel**

Disable interrupts

Enable interrupts

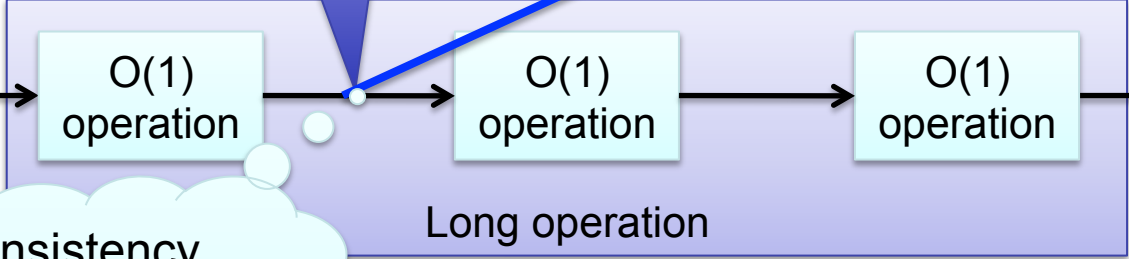
Kernel entry

O(1) operation

Abort & restart later

Kernel exit

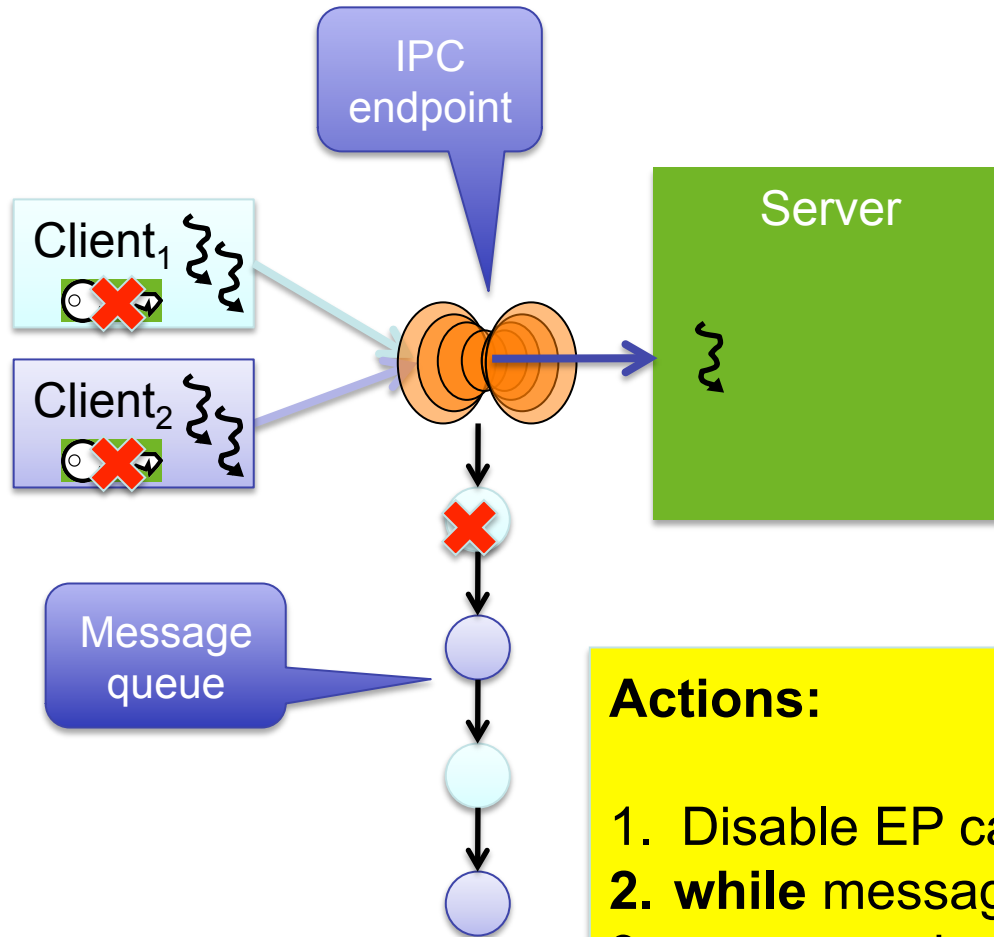
Check pending interrupts



- Consistency
- Restartability
- Progress

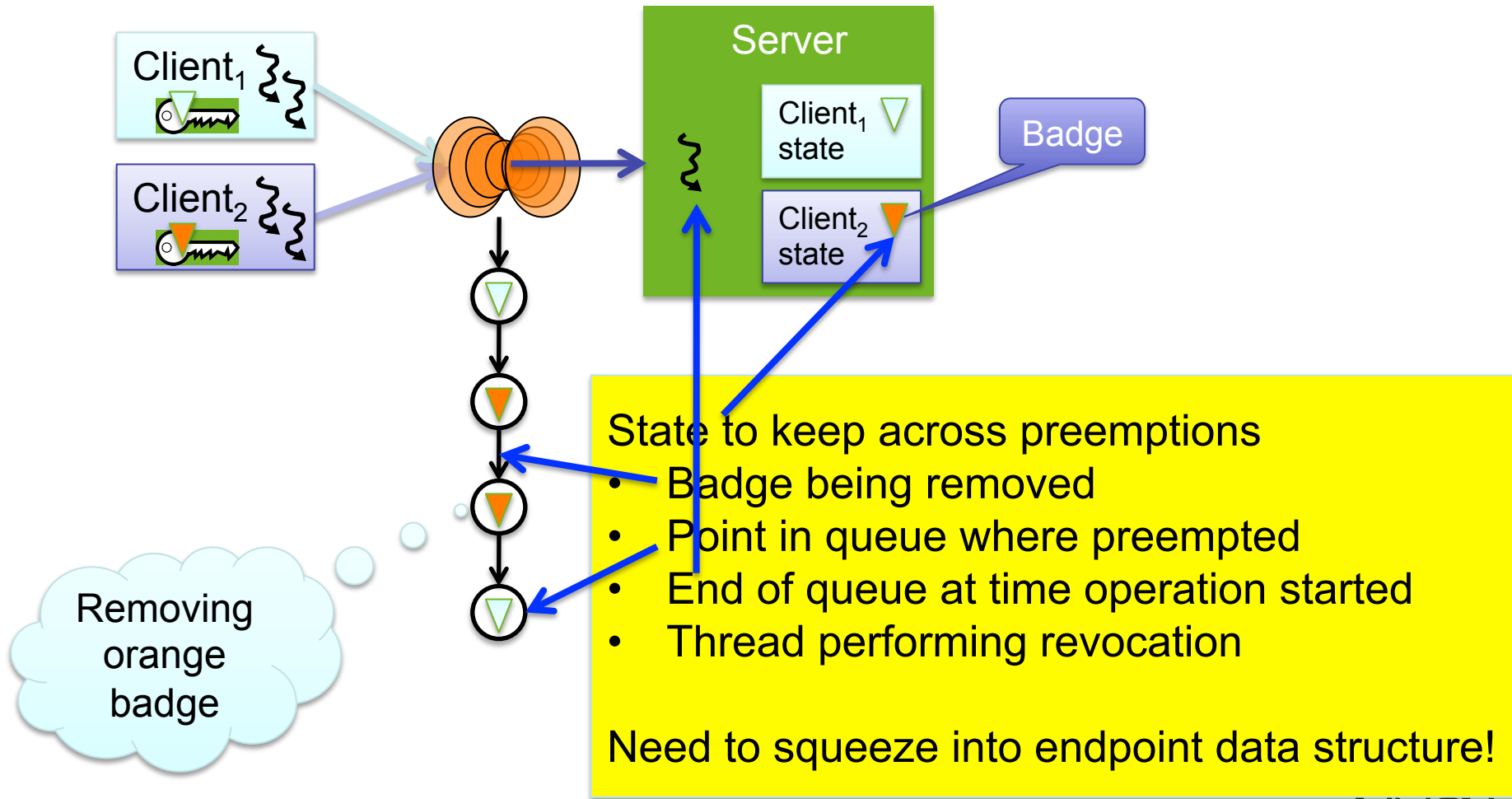


# Example: Destroying IPC Endpoint



- Actions:**
1. Disable EP cap (prevent new messages)
  2. **while** message queue not empty **do**
  3.     remove head of queue (abort message)
  4.     check for pending interrupts
  5. **done**

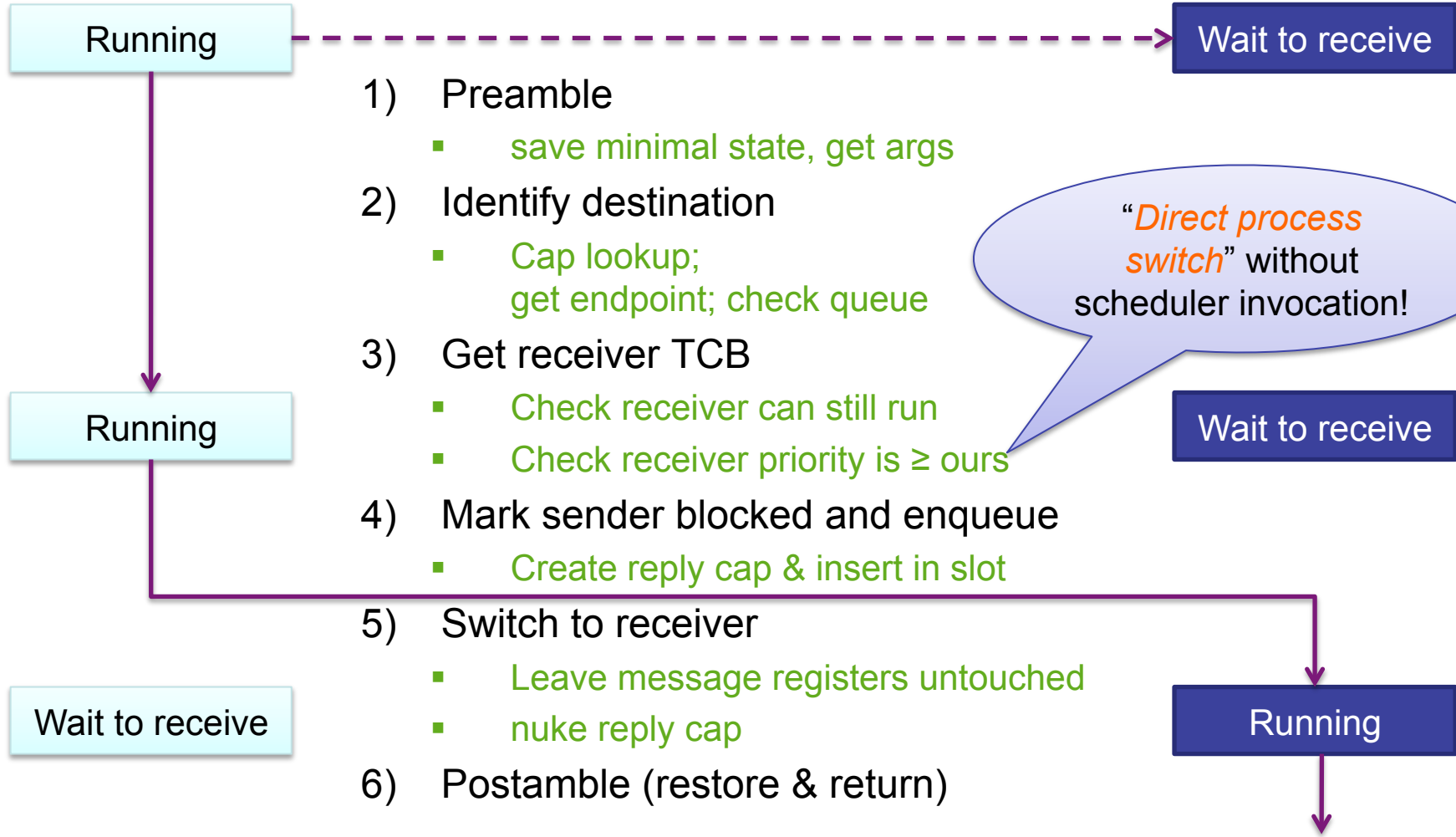
# Difficult Example: Revoking IPC “Badge”



# Synchronous IPC Implementation

185 cycles on ARM11!  
NICTA

Simple send (e.g. as part of RPC-like “call”):



# Fastpath Coding Tricks



```
slow = cap_get_capType(en_c) != cap_endpoint_cap ||  
       !cap_endpoint_cap_get_capCanSend(en_c);  
if (slow) enter_slow_path();
```

Common case: 0

Common case: 1

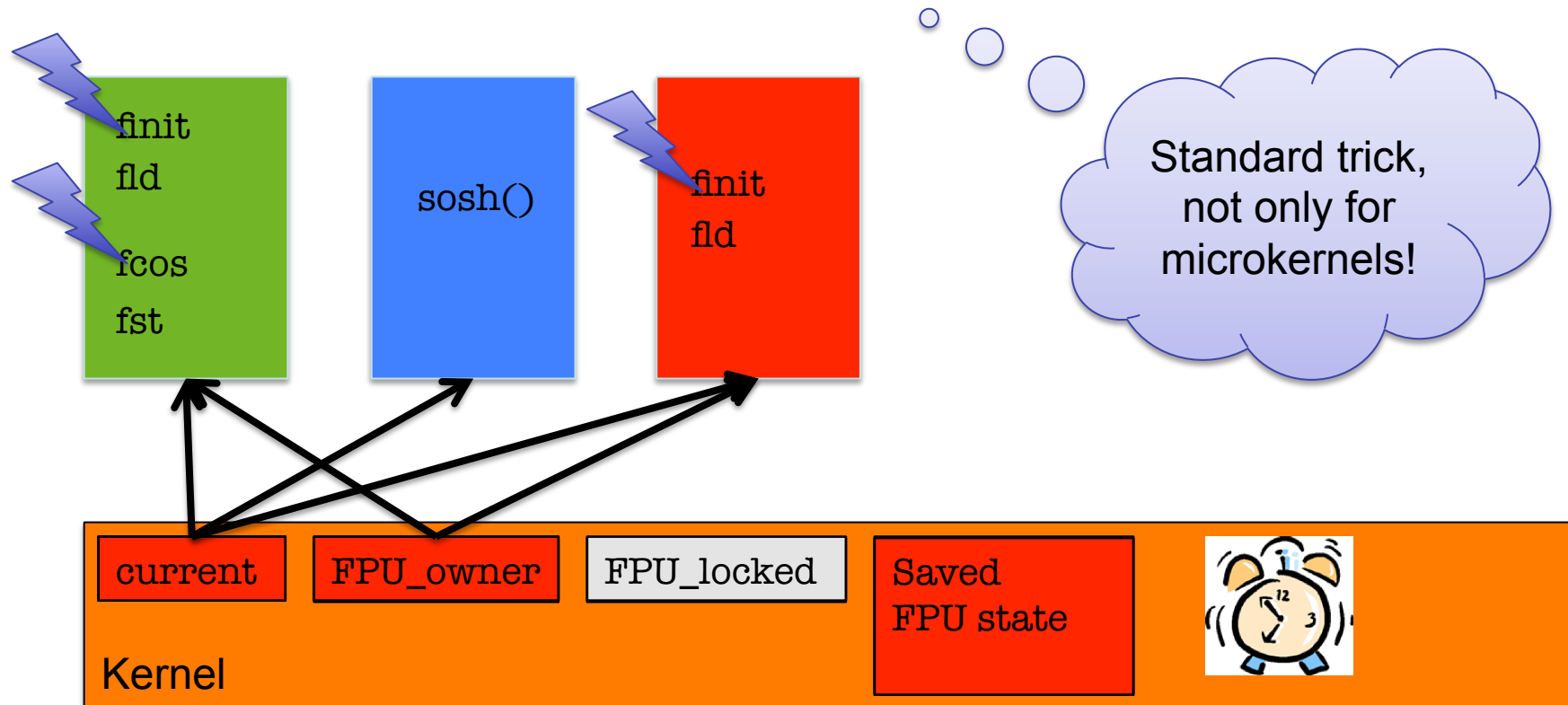
- Reduces branch-prediction footprint
- Avoids mispredicts, stalls & flushes
- Uses ARM instruction predication
- But: increases slow-path latency
  - should be minimal compared to basic slow-path cost



# Lazy FPU Switch



- FPU context tends to be heavyweight
  - eg 512 bytes FPU state on x86
- Only few apps use FPU (and those don't do many syscalls)
  - saving and restoring FPU state on every context switch is wastive!

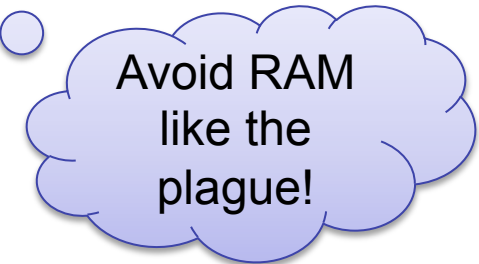


# Other implementation tricks

---



- Cache-friendly data structure layout, especially TCBs
  - data likely used together is on same cache line
  - helps best-case and worst-case performance
- Kernel mappings locked in TLB (using superpages)
  - helps worst-case performance
  - helps establish invariants: page table never walked when in kernel



# Other Lessons Learned from 2<sup>nd</sup> Generation



- Programming languages:
  - original i496 kernel [’95]: all assembler
  - UNSW MIPS and Alpha kernels [’96, ’98]: half assembler, half C
  - Fiasco [TUD ’98], Pistachio [’02]: C++ with assembler “fast path”
  - seL4 [’09], OKL4 [’09]: all C
- Lessons:
  - C++ **ABANDONED**: code bloat, no real benefit
  - Changing calling conventions not worthwhile
    - Conversion cost in library stubs and when entering C in kernel
    - Reduced compiler optimization
  - Assembler **ABANDONED** unnecessary for performance
    - Can write C so compiler will produce near-optimal code
    - C entry from assembler cheap if calling conventions maintained
    - seL4 performance with C-only passthrough as good as other L4 kernels [Blackham & Heiser ’12]

# L4 Design and Implementation



## Implement. Tricks [SOSP'93]

- ~~Process kernel~~
- ~~Virtual TCB array~~
- ~~Lazy scheduling~~
- Direct process switch
- Non-preemptible
- ~~Non-portable~~
- ~~Non-standard calling convention~~
- ~~Assembler~~

## Design Decisions [SOSP'95]

- Synchronous IPC
- Rich message structure, ~~arbitrary out-of-line messages~~
- Zero-copy register messages
- User-mode page-fault handlers
- ~~Threads as IPC destinations~~
- ~~IPC timeouts~~
- ~~Hierarchical IPC control~~
- User-mode device drivers
- ~~Process hierarchy~~
- ~~Recursive address-space construction~~

# seL4 Design Principles

---



- Fully delegatable access control
- All resource management is subject to user-defined policies
  - Applies to kernel resources too!
- Suitable for *formal verification*
  - Requires small size, avoid complex constructs
- Performance on par with best-performing L4 kernels
  - Prerequisite for real-world deployment!
- Suitability for real-time use
  - Only partially achieved to date ☹
    - on-going work...

# (Informal) Requirements for Formal Verification

---



- Verification scales poorly  $\Rightarrow$  small size (LOC and API)
- Conceptual complexity hurts  $\Rightarrow$  KISS
- Global invariants are expensive  $\Rightarrow$  KISS
- Concurrency difficult to reason about  $\Rightarrow$  single-threaded kernel

Largely in line with traditional L4 approach!

# seL4 Fundamental Abstractions

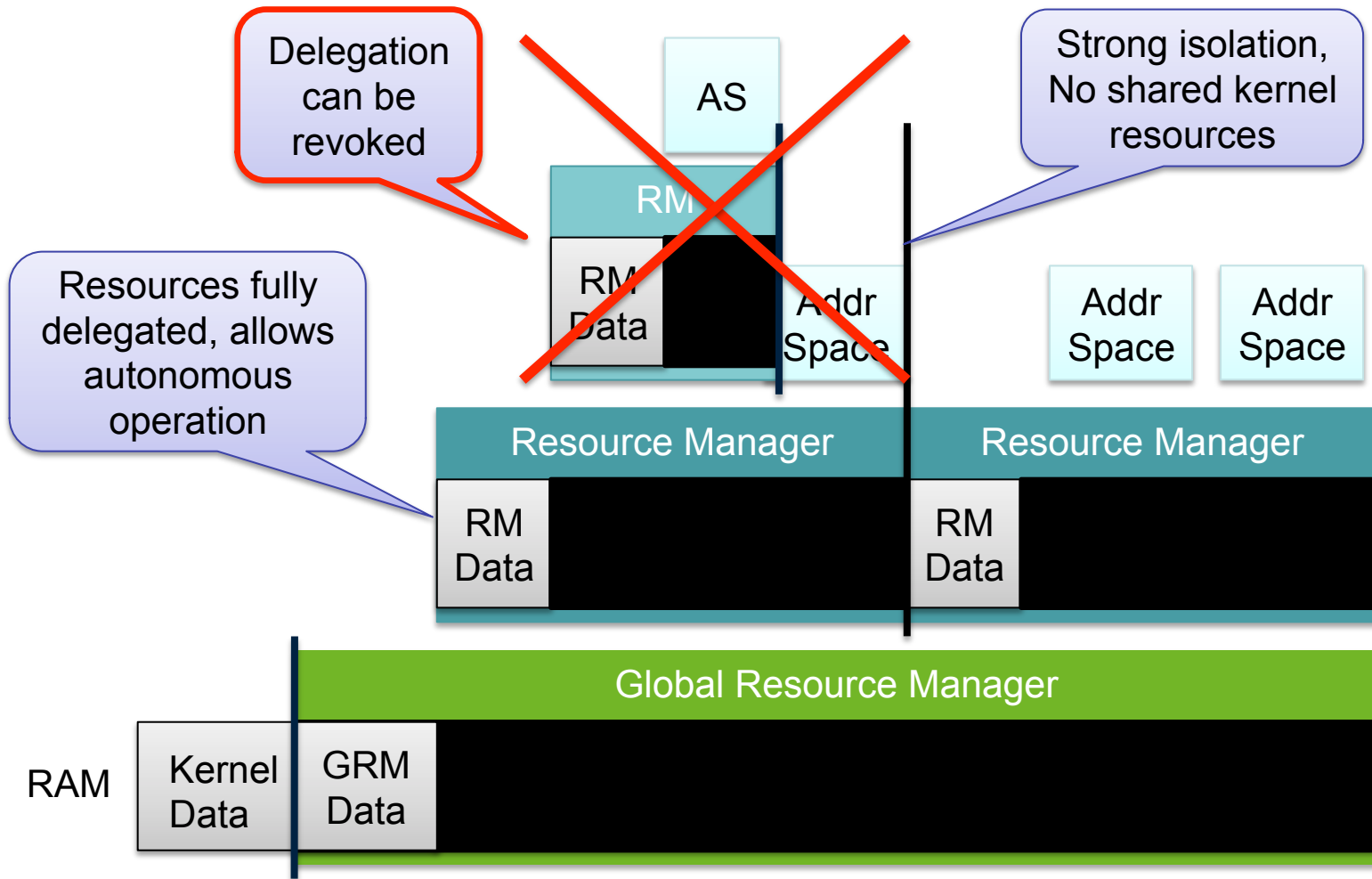


- Capabilities as opaque names and access tokens
  - All kernel operations are cap invokations (except Yield())
- IPC:
  - Synchronous (blocking) message passing plus asynchronous notification
  - Endpoint objects implemented as message queues
    - Send: get receiver TCB from endpoint or enqueue self
    - Receive: obtain sender's TCB from endpoint or enqueue self
- Other APIs:
  - Send()/Receive() to/from virtual kernel endpoint
  - Can interpose operations by substituting actual endpoint
- Fully user-controlled memory management

seL4's main conceptual novelty!

A light blue thought bubble with a white outline and a drop shadow. It is connected to the text "Fully user-controlled memory management" by three small circles of increasing size.

# Remember: seL4 User-Level Memory Management





# Remaining Conceptual Issues in seL4

---

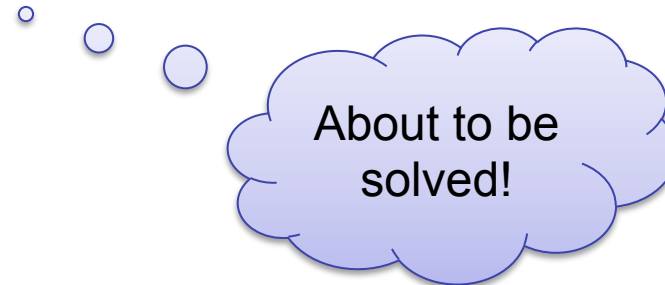


## IPC & Tread Model:

- Is the “mostly synchronous + a bit of async” model appropriate?
  - forces kernel scheduling of user activities
  - forces multi-threaded userland

## Time management:

- Present scheduling model is ad-hoc and insufficient
  - fixed-prio round-robin forces policy
  - not sufficient for some classes of real-time systems (time triggered)
  - no real support for hierarchical real-time scheduling
  - lack of an elegant resource management model for time



# Lessons From 20 Years of L4

---



- Minimality is excellent driver of design decisions
  - L4 kernels have become simpler over time
  - Policy-mechanism separation (user-mode page-fault handlers)
  - Device drivers really belong to user level
  - Minimality is key enabler for formal verification!
- IPC speed still matters
  - But not everywhere, premature optimisation is wasteful
  - Compilers have got so much better
  - Verification does not compromise performance
  - Verification invariants can help improve speed! [Shi, OOPSLA'13]
- Capabilities are the way to go

- Details changed, but principles remained
- Microkernels rock! (If done right!)